



# A Mean-Field Model for Multiple TCP Connections through a Buffer Implementing RED

François Baccelli, David R. McDonald, Julien Reynier

## ► To cite this version:

François Baccelli, David R. McDonald, Julien Reynier. A Mean-Field Model for Multiple TCP Connections through a Buffer Implementing RED. [Research Report] RR-4449, INRIA. 2002. inria-00072139

**HAL Id: inria-00072139**

**<https://inria.hal.science/inria-00072139>**

Submitted on 23 May 2006

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



INSTITUT NATIONAL DE RECHERCHE EN INFORMATIQUE ET EN AUTOMATIQUE

# *A Mean-Field Model for Multiple TCP Connections through a Buffer Implementing RED*

François Baccelli — David R. McDonald — Julien Reynier

N° 4449

April 2002

THÈME 1

A large blue rectangle occupies the lower half of the page. Overlaid on the left side of this rectangle is a large, light grey stylized letter 'R'. To the right of the 'R', the words 'apport de recherche' are written in a white, italicized serif font. A horizontal grey brushstroke underline is positioned beneath the text.

*apport  
de recherche*





## A Mean-Field Model for Multiple TCP Connections through a Buffer Implementing RED

François Baccelli <sup>\*</sup>, David R. McDonald <sup>†</sup>, Julien Reynier <sup>‡</sup>

Thème 1 — Réseaux et systèmes  
Projet TREC

Rapport de recherche n° 4449 — April 2002 — 19 pages

**Abstract:** Active queue management schemes like RED (Random Early Detection) have been suggested when multiple TCP sessions are multiplexed through a bottleneck buffer. The idea is to detect congestion before the buffer overflows and packets are lost. When the queue length reaches a certain threshold RED schemes drop/mark incoming packets with a probability that increases as the queue size increases. The objectives are an equitable distribution of packet loss, reduced delay and delay variation and improved network utilization.

Here we model multiple connections maintained in the congestion avoidance regime by the RED mechanism. The window sizes of each TCP session evolve like independent dynamical systems coupled by the queue length at the buffer. We introduce a mean-field approximation to one such RED system as the number of flows tends to infinity. The deterministic limiting system is described by a transport equation. The numerical solution of the limiting system is found to provide a good description of the evolution of the distribution of the window sizes, the average queue size, the average loss rate per connection and the total throughput. TCP with RED or tail-drop may exhibit limit cycles and this causes unnecessary packet delay variation and variable loss rates. The root cause of these limit cycles is the hysteresis due to the round trip time delay in reacting to a packet loss.

**Key-words:** TCP, RED, mean-field, dynamical system, transport equation.

<sup>\*</sup> INRIA-ENS, ENS, 45 rue d'Ulm 75005 Paris, France {Francois.Baccelli@ens.fr}

<sup>†</sup> Department of Mathematics, University of Ottawa, Canada {dmdsg@mathstat.uottawa.ca}

<sup>‡</sup> INRIA-ENS, ENS, 45 rue d'Ulm 75005 Paris, France {Julien.Reynier@ens.fr}

## **Un Modèle de Champ Moyen pour Plusieurs Connexions TCP se Partageant une File d'Attente RED**

**Résumé :** Des schémas de gestion active tels que RED (Random Early Detection) ont été proposés pour les files d'attente en amont d'une ressource partagée entre un ensemble de sessions contrôlées par TCP. Ces schémas sont fondés sur la détection du franchissement de certains seuils de congestion qui anticipent la perte de paquets due au débordement de la file d'attente. Lorsque le remplissage de cette dernière dépasse un seuil, RED détruit ou marque tout nouveau paquet arrivant dans la file avec une probabilité qui augmente avec le niveau du seuil. L'objectif est de distribuer de manière équitable les pertes de paquets, de réduire les délais et leurs variations et d'améliorer l'utilisation du réseau.

Dans cet article, nous étudions l'effet de RED sur un ensemble de sessions TCP, toutes dans la phase d'évitement de congestion. Les fenêtres de chacune des sessions évoluent comme des systèmes dynamiques qui ne sont couplés que par le niveau de la file d'attente partagée. Nous introduisons une approximation de champ moyen en faisant tendre le nombre des sessions vers l'infini. Le système limite est régi par une équation de transport déterministe. La solution numérique de ce système limite fournit une bonne description de l'évolution de la distribution des fenêtres, ainsi que de celle la moyenne de la file d'attente, du taux de perte moyen par session et du débit total. Que ce soit avec RED ou avec les pertes par débordement, TCP est sujet à des cycles limites qui peuvent causer des oscillations des délais et des taux de pertes. Ces cycles limites sont dûs à un phénomène d'hystérésis: pour toute session TCP, le délai de réaction à une perte de paquet est en effet de l'ordre d'un RTT.

**Mots-clés :** TCP, RED, champ moyen, système dynamique, équation de transport.

## 1 Introduction

Active Queue Management and in particular Random Early Detection (RED) schemes have been proposed to improve on the basic tail drop mechanism. In RED, an arriving packet is killed with a probability which increases with the queue size. The original RED scheme (see ([6])) proposed a linear increase in loss rate from 0 at queue size  $Q_{min}$  to  $p_{max}$  at  $Q_{max}$  and 1 for a queue size greater than  $Q_{max}$ . Other RED schemes have been suggested in [2] and [9].

Early detection of incipient congestion causes some connections to reduce their transmission rate (due to a packet loss) long before buffer overflow, one hopes there will be a high utilization of the network with an equitable distribution of packet loss since those connections with large windows have a higher chance of incurring losses.

In the present paper, we study the interaction of a large number  $N$  of TCP/IP connections controlled by TCP Reno, which are all routed through a bottleneck queue in a router implementing RED. An important feature is that there is a delay of one round trip time between the time the packet is killed and the time when the buffer receives the reduced rate. During this delay time packets continue to arrive at the old rate. This has been seen ([2]) to result in hysteresis effects that include buffer oscillations. These oscillations appear as limit cycles in our theoretical model.

Section 2 contains the model notation and the main assumptions. In Subsection 3.1 we construct the model describing the evolution of the window sizes and the queue size. In Subsection 3.2 we reformulate the description of the window sizes as a random measure or random histogram which is coupled to the queue at the buffer. In Subsection 3.3 we show that this model converges, when the number of sessions  $N$  tends to infinity, to a deterministic transport equation (3.10) and (3.11) describing the evolution of the deterministic histogram of window sizes coupled with a deterministic fluid queue size.

Such mean field limits have long been discussed in the statistical physics literature (see [4]) and recently in the queueing literature [8]. Our mean field limit is a transport equation which has also appeared in the context of cell biology where one wishes to model the distribution of the sizes of a population of cells which have linear growth followed by cell mitosis which occurs with a probability depending on the size of the cell (see [5]).

Our model is identical to that in [9] (which only discusses RED) but we keep the histogram of the window sizes as our state descriptor while [9] only keeps the mean of this histogram. Taking the mean window size collapses our equations to Equation (1) in [9] under the hypothesis that the window distributions at any time  $t$  and one round trip time before  $t$  are *uncorrelated* (see (3.12)). A similar model describing the evolution of the mean window size is given in [10] (see Equation (1) there). Also see ([11]). Our model and our analysis of the mean field limit of the histogram of the window size is similar to that in [1] except the later does not include acknowledgement delay and does not discuss RED. It should be emphasized that our model handles both tail-drop or RED (and a combination of both) in the same framework. Our model can easily be adapted to describe Explicit Congestion Notification (ECN) since it can handle high marking rates where one would expect a strong correlation between present and recent window sizes.

To approximate the performance of a given (finite speed) router with a given number of connections using our mean field model we simply define  $L$  to be the router's link rate divided by the number of active connections. With this normalization our mean field model provides a good fit to the trace of queue sizes obtained from Opnet simulations. For some parameter choices the mean field system (and the Opnet trace) stabilizes as is shown in Figures 3 and 2 respectively in Section 5 but by increasing the round trip time delays the mean field system (and the Opnet trace) become unstable as in Figures 5. This phenomenon was investigated in [9] and it results from an amplification of perturbations due

to feedback delay. When the system is stable we can give a closed form expression for the window size histogram (see (4.18)). The window distribution provides QoS results like the proportion of time a connection has a throughput less than any particular value. It also provides a means of calculating the proportion of connections in timeout due to successive packet losses or due to a packet loss when the congestion window is less than four packets (see (3) for a closed form expression).

A rigorous proof of the convergence and the existence of the mean field limit along the lines of [8] is beyond the scope of this paper. We only give a plausibility argument for the resulting equations.

## 2 Notation and Assumptions

We study the interaction of  $N$  TCP/IP connections controlled by TCP Reno, which are routed through a bottleneck queue in a router. Hence each of the connections implements a window flow control which limits the number of packets from this connection allowed into the network during one Round Trip Time (RTT). The link rate of the router is  $NL$  packets per second.

We assume the packets from all active connections join the queue  $Q(t)$  at the bottleneck buffer and we denote by  $Q_N(t)$  the average queue per flow so the length of the queue at time  $t$  is  $Q(t) = NQ_N(t)$ . We assume the scheduling to be FIFO.

We imagine the source writes its current window size and the current RTT in each packet it sends and we define

- $W_n(t)$  the window size written in a packet from connection  $n$  arriving at the server at time  $t$ ;
- $R_n(t)$  the RTT written in a packet from connection  $n$  arriving at the server at time  $t$  (this RTT is the sum of the propagation delay plus the queueing delay in the router).

Let  $\mathbf{W}(t) = (W_1(t), \dots, W_N(t))$  represent the state (the window sizes) of all connections at time  $t$ . Throughout the paper, we will approximate the real system by saying the throughput at time  $t$  is the window size at time  $t$  divided by the RTT at this time.

Under TCP Reno, established connections execute congestion avoidance where the window size of each connection increases by one packet each time a packet makes a round trip, i.e. each  $R_n$  as long as no losses or timeouts occur. During this phase the rate the window of connection  $n$  increases is approximately  $1/R_n$  packets per second. The only thing restraining the growth of transmission rates is a loss or timeout.

In the present paper, we will neglect slow start and a negotiated maximum window size although the mean field method can easily be adapted to take these features into account by enlarging the state space to be a multivariate histogram describing the joint distribution of the window size with, say, the slow start thresholds. We will also assume there are no transmission losses. Hence the only losses or explicit congestion notifications (ECN) are generated by active buffer management or by tail-drop. When a loss or ECN occurs the window is reduced by half.

We will assume the buffer has size  $B$  packets and that once this buffer space is exhausted arriving packets are dropped. Such tail-drops come in addition to the RED mechanism. Here we take the drop probability of RED (of an incoming packet before being processed) to be a function  $\mathcal{F}(Q(t))$  which is zero for  $Q(t)$  below  $Q_{min}$  but rises to  $p_{max}$  at  $Q_{max}$  and further to 1 when  $Q(t)$  reaches  $B$ . If the number of active connections is  $N$  we can reformulate this drop probability as  $F(Q_N(t))$ , where  $F$  is a distribution function which is zero below  $q_{min} = Q_{min}/N$  but rises to  $p_{max}$  at  $q_{max} = Q_{max}/N$  and further to 1 when  $Q_N(t)$  reaches  $B$  where  $B = B/N$ . Of course the tail-drop scheme can be considered as the limiting case when  $F(q) = 0$  for  $q < B$  and  $F(q) = 1$  for  $q \geq B$ .

As we shall see, the model takes into account that delay of one round trip time between the time the packet is killed and the time when the buffer receives the reduced rate, and leads to oscillations for some values of the parameters. The best  $F$  would minimize the dispersion of the window sizes and eliminate oscillations of the queue length (thus reducing packet delay variation).

### 3 The $N$ -particle system and mean-field limit

#### 3.1 The $N$ -Particle Markov process

We assume window reductions at connection  $n$  occur according to a Poisson process with stochastic intensity

$$\frac{W_n(t - R_n(t))}{R_n(t)} F(Q_N(t - R_n(t)))$$

(we can assume  $W_n(t) = 0$  for  $t < 0$ ). This makes the loss rate proportional to the transmission rate one RTT in the past multiplied by the RED loss rate one RTT in the past and thus imitates reality. We could describe the RED mechanism of using a moving exponential weighted average of past queue sizes to determine the drop rate. In this case the stochastic intensity would be given by

$$\frac{W_n(t - R_n(t))}{R_n(t)} F \left( \int_{-\infty}^{t - R_n(t)} Q_N(s) \exp(-\beta(t - R_n(t) - s)) ds \right)$$

where  $\exp(-\beta)$  is the exponential averaging coefficient.

Let  $\{N_n(t); n = 1, \dots, N\}$  be  $N$  independent Poisson processes with intensity 1 and let

$$\Lambda_n(t) = \int_0^t \frac{W_n(s - R_n(s))}{R_n(s)} F(Q_N(s - R_n(s))) ds$$

be the stochastic intensity for the Poisson point process of losses of connection  $n$ . Hence the losses of connection  $n$  occur according to the time changed Poisson process  $N_n(\Lambda_n(t))$ .

When no loss occurs the window size increases linearly at rate  $1/R_n(t)$  so the window size is increased by approximately  $1/W_n(t - R_n(t))$  each time an acknowledgement returns to the source. But when a packet was lost at time  $t - R_n(t)$ , the source does not increase the window size this amount, and in addition it cuts the current size by half. Hence the evolution of the window size is described by the following stochastic differential equation:

$$dW_n(t) = \frac{1}{R_n(t)} dt - \left( \frac{1}{W_n((t - R_n(t))^-)} + W_n(t^-)/2 \right) dN_n(\Lambda_n(t)), \quad (3.1)$$

with  $W_n(0) = w_n(0)$ ,  $n = 1, \dots, N$  specified. Note that  $t^-$  means the left limiting value at  $t$ . Also note that we would incorporate a negotiated maximum queue size  $w_{max}$  by modifying the above drift term to  $\chi\{W_n(t) \leq w_{max}\} \cdot 1/R_n(t)$  where  $\chi$  denotes the indicator function.

As a rough approximation, because of the FIFO assumption,  $R_n$  should satisfy  $R_n(t) = T_n + Q_N(t - R_n(t))/L$ , where  $T_n$  is the propagation delay from source  $n$  to the destination and back. Note that  $W_n(t^-)$  is completely determined by  $\mathcal{F}_n(t)$ , the past one RTT ago generated by  $\{W_n(s - R_n(s)), N_n(\Lambda_n(s)), 0 \leq s \leq t^-\}$ .

It will be easier to approximate the above dynamical system by a fluid model, where the queue, the windows and the thresholds to evolve as a differential system. We assume packets have equal mean sizes of 1 data unit. When there are no losses, the rate at which source  $n$  pours fluid into the buffer



is  $W_n(t)/R_n(t)$ . A loss at time  $t - R_n(t)$  means the source stops sending packets until  $W_n(t)/2$  packets are acknowledged, i.e. until the window size has been reduced by half. As far as the queue is concerned, it sees an throughput of  $W_n(t)/R_n(t)$  for roughly half an  $RTT$  since packets in the system continue to arrive at the old rate. This is followed by a zero throughput for the remaining half of the  $RTT$ . Hence the average throughput for an  $RTT$  following a loss is  $W_n(t)/(2R_n(t))$ . Since the window size is halved when a loss is detected, according to our convention, the throughput over the  $RTT$  following a loss is  $W_n(t)/(2R_n(t))$ , i.e. equal to that of the real system. Hence the rate of change of the fluid buffer is given by

$$\begin{aligned} N \frac{dQ_N(t)}{dt} &= \sum_{n=1}^N \frac{W_n(t)}{R_n(t)} (1 - F(Q_N(t))) - NL \\ &\quad + \left( \sum_{n=1}^N \frac{W_n(t)}{R_n(t)} (1 - F(Q_N(t))) - NL \right)^- \chi\{Q_N(t) = 0\} \end{aligned}$$

since the proportion  $F(Q_N(t))$  of the fluid is lost. The second term prevents the queue size from becoming negative. In effect the queue can stick at 0 until a sufficient number of connections increase their window size. Dividing by  $N$  gives

$$\begin{aligned} \frac{dQ_N(t)}{dt} &= \frac{1}{N} \sum_{n=1}^N \frac{W_n(t)}{R_n(t)} (1 - F(Q_N(t))) - L \\ &\quad + \left( \frac{1}{N} \sum_{n=1}^N \frac{W_n(t)}{R_n(t)} (1 - F(Q_N(t))) - L \right)^- \chi\{Q_N(t) = 0\}, \end{aligned} \quad (3.2)$$

with  $Q_N(0) = q(0)$ . Note that under ECN there is no loss associated with marking so the term  $(1 - F(Q_N(t)))$  in (3.2) disappears.

### 3.2 Reformulation in terms of a measure-valued process

At this point we make the simplifying assumption that all connections have a constant transmission time  $T_n = T$ .

Consequently  $R(t) = T + Q_N(t - R(t))/L$ . In order to study the limiting behavior of the system (3.1), (3.2) as the number of connections  $N$  goes to infinity, we reformulate the system in terms of its empirical process (see Dawson [4]). For any Borel set  $A \subset \hat{S}$  define

$$M_N(t, A) := \frac{1}{N} \sum_{n=1}^N \chi_A(W_n(t)) \quad (3.3)$$

to be the associated probability-measure-valued process.

The process  $M_N(t) \equiv M_N(t, \cdot)$  belongs to the state space  $M_1(\mathbb{R}^+)$ , the set of probability measures on  $\mathbb{R}^+ = [0, \infty)$  furnished with the topology of weak convergence. Given an initial distribution  $\mu(A) = \frac{1}{N} \sum_{n=1}^N \chi_A(w_n(0))$  in  $M_1(\mathbb{R}^+)$  and an initial value of  $Q_N(0) = q(0)$  specifies the canonical process  $(M, Q)$  with marginals  $(M_N(t, \cdot), Q_N(t))$  on the set of trajectories  $\Omega = C([0, \infty), M_1(\mathbb{R}^+) \times \mathbb{R}^+)$ , the space of continuous functions from  $[0, \infty)$  into  $M_1(\mathbb{R}^+) \times \mathbb{R}^+$ . Let  $P_{\mu, q(0)}$  denote the induced probability measure on  $\Omega$ . We shall also need a joint measure-valued process

$$M_N(s, B; t, A) := \frac{1}{N} \sum_{n=1}^N \chi_B(W_n(s)) \chi_A(W_n(t)).$$

It will be clear from context if  $M_N$  denotes the joint or marginal process.

Since the dynamical systems specified by (3.1), (3.2) are exchangeable in the  $W_n(t)$  (we can relabel the connections without changing the evolution of the system) it follows that (3.1), (3.2) can be reformulated in terms of  $(M_N(t), Q_N(t))$ . The dynamics of  $M_N(t)$  are described through a set of equations satisfied by the scalar product

$$\langle g, M_N(t) \rangle = \int_0^\infty g(w) M_N(t, dw)$$

where  $g \in \mathcal{G}$  and  $\mathcal{G} = \{g \in C_b^1(\mathbb{R}^+) : g(0) = 0\}$  with  $C_b^1(\mathbb{R}^+)$  the space of bounded functions with bounded derivatives.  $\mathcal{G}$  is chosen to avoid singular behaviour associated with connections disappearing from the system.

Reformulating (3.2) we get

$$\begin{aligned} & Q_N(t) - Q_N(0) \\ &= \int_0^t \left[ \langle w, M_N(s) \rangle \frac{(1 - F(Q_N(s)))}{R(s)} - L \right. \\ & \quad \left. + \left( \langle w, M_N(s) \rangle \frac{1 - F(Q_N(s))}{R(s)} - L \right)^- \chi\{Q_N(s) = 0\} \right] ds \\ &= \int_0^t \left[ \langle w, M_N(s) \rangle \frac{1}{R(s)} (ds - dK_N(s)) - Lds \right. \\ & \quad \left. + \left( \langle w, M_N(s) \rangle \frac{1}{R(s)} (ds - dK_N(s)) - Lds \right)^- \chi\{Q_N(s) = 0\} \right], \end{aligned} \quad (3.4)$$

where  $K_N(t) = \int_0^t F(Q_N(s))ds$  is the cumulative loss or kill rate.

We can also reformulate (3.1):

$$\begin{aligned} & \langle g, M_N(t) \rangle - \langle g, M_N(0) \rangle \\ &= \frac{1}{N} \sum_{n=1}^N \int_0^t \left[ \frac{dg}{dw}(W_n(s)) \frac{1}{R(s)} ds \right. \\ & \quad \left. + \left( -\frac{dg}{dw}(W_n(s^-)) \frac{1}{W_n((s - R(s))^-)} + g(W_n(s^-)/2) - g(W_n(s^-)) \right) dN_n(\Lambda_n(s)) \right] \\ &= \frac{1}{N} \sum_{n=1}^N \int_0^t \left[ \frac{dg}{dw}(W_n(s)) \frac{1}{R(s)} ds \right. \\ & \quad \left. + \left( -\frac{dg}{dw}(W_n(s)) \frac{1}{W_n(s - R(s))} + g(W_n(s)/2) - g(W_n(s)) \right) \right. \\ & \quad \left. \cdot \frac{W_n(s - R(s))}{R(s)} F(Q_N(s - R(s))) ds \right] + \mathcal{E}_N(t) \end{aligned}$$

where

$$\mathcal{E}_N(t) = \sum_{n=1}^N \int_0^t \left( -\frac{dg}{dw}(W_n(s^-)) \frac{1}{W_n((s - R(s))^-)} + g(W_n(s^-)/2) - g(W_n(s^-)) \right) dZ_n(\Lambda_n(s))$$

and

$$Z_n(t) - Z_n(0) := \int_0^t \left( dN_n(\Lambda_n(s)) - \frac{W_n(s - R(s))}{R(s)} F(Q_N(s - R(s))) ds \right).$$

Hence,

$$\begin{aligned} & \langle g, M_N(t) \rangle - \langle g, M_N(0) \rangle \\ &= \int_0^t \left[ \frac{1}{R(s)} (1 - F(Q_N(s - R(s)))) \left\langle \frac{dg(w)}{dw}, M_N(s) \right\rangle \right. \end{aligned} \quad (3.5)$$

$$\begin{aligned} & \left. + \langle (g(w/2) - g(w))v, M_N(s - R(s), dv; s, dw) \rangle \frac{1}{R(s)} F(Q_N(s - R(s))) \right] ds \\ &+ \mathcal{E}_N(t) \\ &= \int_0^t \left[ \frac{1}{R(s)} (1 - F(Q_N(s - R(s)))) \left\langle \frac{dg(w)}{dw}, M_N(s) \right\rangle ds \right. \\ & \left. + \langle (g(w/2) - g(w))v, M_N(s - R(s), dv; s, dw) \rangle \frac{1}{R(s)} dK_N(s - R(s)) \right] \quad (3.6) \\ &+ \mathcal{E}_N(t). \end{aligned}$$

### 3.3 The mean-field evolution equations

As the number of connections  $N$  becomes large a remarkable simplification occurs essentially because of the law of large numbers. The error term  $\mathcal{E}_N(t)$  is a martingale with mean value 0 whose supremum over any bounded interval of time converges to 0 in probability. This leaves behind a deterministic system. Hence in the limit the histogram of the window sizes becomes deterministic as does the queue size and the resulting deterministic mean field system is described in the following result.

**Theorem 1** *Suppose that as  $N \rightarrow \infty$ ,  $\mu_N = M_N(0)$  converges weakly to some  $\mu(0) \in M_1(\mathbb{R}^+)$  and  $Q_N(0)$  converges to  $q(0)$ . Then  $M_N(t) \rightarrow \mu(t)$ ,  $Q_N(t) \rightarrow q(t)$  and  $K_N(t) \rightarrow K(t)$  where  $\mu(t)$ ,  $q(t)$  and  $K(t)$  are continuous functions of  $t \in \mathbb{R}^+$  into  $M_1(\mathbb{R}^+)$  and  $\mathbb{R}^+$  respectively and  $K(t) = \int_0^t k(s)ds$  so  $F(Q_N(t)) \rightarrow k(t)$  at points of continuity of  $F$ . Moreover, for any function  $g \in \mathcal{G}$ ,*

$$\begin{aligned} & \langle g, \mu(t) \rangle - \langle g, \mu(0) \rangle \\ &= \int_0^t \left[ \frac{1}{r(s)} (1 - k(s - r(s))) \left\langle \frac{dg(w)}{dw}, \mu(s) \right\rangle \right. \end{aligned} \quad (3.7)$$

$$\begin{aligned} & \left. + \langle (g(w/2) - g(w))v, \mu(s - r(s), dv; s, dw) \rangle \frac{1}{r(s)} k(s - r(s)) \right] ds \\ &= \int_0^t \left[ \frac{1}{r(s)} (1 - k(s - r(s))) \left\langle \frac{dg(w)}{dw}, \mu(s, dw) \right\rangle \right. \quad (3.8) \\ & \left. + \langle (g(w)v, \mu(s - r(s), dv; s, d(2w)) - \mu(s - r(s), dv; s, dw)) \rangle \frac{1}{r(s)} k(s - r(s)) \right] ds \end{aligned}$$

and

$$\begin{aligned} & q(t) - q(0) \quad (3.9) \\ &= \int_0^t \left[ \langle w, \mu(s) \rangle \frac{1}{r(s)} (1 - k(s)) - L + \left( \langle w, \mu(s) \rangle \frac{1}{r(s)} (1 - k(0)) - L \right)^- \chi\{q(s) = 0\} \right] ds \end{aligned}$$

where  $r(t) = T + q(t - r(t))/L$ .

Note that these equations do permit solutions where  $q(t)$  reaches and sticks to the maximum boundary  $B$  or  $q_{max}$  where  $F$  has a discontinuity. Tail-drop is the most obvious example. In this case  $Q_N(t)$  jitters at and below this boundary and the loss rate  $F(Q_N(t))$  jitters between the values 0 and 1. Since we have weak convergence of the cumulative loss rate  $\int_0^t F(Q_N(s))ds$  to the deterministic limiting cumulative loss rate  $K_N(t) = \int_0^t k(s)ds$  we may consider  $k(t)$  is effective loss rate at time  $t$ . These equations also allow for the case when  $F(0+) \neq 0$  and  $q(t)$  reaches and sticks to zero. Again the queue will jitter and the effective loss rate is  $k(0)$ .

$\mu(s, dv; t, dw)$  and  $q(t)$  satisfy equations (3.8) and (3.9) but these equations do not determine  $\mu(s, dv; t, dw)$  and  $q(t)$ . We can consider a larger state space with states representing the trajectory of the window histograms from one round trip time before  $t$  up to time  $t$ . The resulting system is Markovian and  $\mu(t - r(t), dv; t, dw)$  is given by the joint marginal distribution at times  $t - r(t)$  and  $t$ . This enlarged system is useful for the proof of convergence but useless for practical purposes.

Finally note that the proof of the above theorem is incomplete so despite its plausibility it is really still a conjecture.

The fluctuations observed in the Opnet simulations with increasing  $N$  as in Figures 2 ( $N = 200$ ), Figure 4 ( $N = 400$  and  $N = 800$ ) diminish to zero in the mean field limit as  $N$  increases. The investigation of these fluctuations rescaled by a factor of  $\sqrt{N}$  as in [4] has not been done.

We can make simplifying approximations which give some insight into the solution. We make the approximation

$$E(W_n(s - R(s)) | W_n(s), Q_N(s - R(s))) \approx W_n(s).$$

This is equivalent to assuming

$$\int_v \mu(s - r(s), v; s, dw) v dv = w \mu(s, dw).$$

This is clearly inaccurate if the loss rate is very small for then the window size one RTT ago would be one less than it is now. Moreover, when the loss rate is moderate, bigger windows are more likely to have been twice as big on RTT ago and have suffered a loss since. Nevertheless, if the loss rate is low and there is a stable fixed point for the system (3.8) and (3.9) then close to the fixed point the expected value of the window size one RTT ago would be the current window size. Consequently, the above approximation will yield a system with the same fixed point. This approximation is given by the following simplified system.

**Corollary 1** *If the initial distribution  $\mu(t, dw)$  has a continuous density then  $\mu(t, dw)$  has a continuous density  $p(t, w)$  differentiable in  $t$  which approximately satisfies the following equations :*

$$\begin{aligned} \frac{\partial p(t, w)}{\partial t} &= \left( p(t, 2w) \frac{2w}{r(t)} - p(t, w) \frac{w}{r(t)} \right) k(t - r(t)) \\ &\quad - \frac{1}{r(t)} (1 - k(t - r(t))) \frac{\partial p(t, w)}{\partial w} \end{aligned} \quad (3.10)$$

and

$$\begin{aligned} \frac{dq(t)}{dt} &= \int_w \frac{w}{r(t)} p(t, w) dw (1 - k(t)) - L \\ &\quad - \left( \int_w \frac{w}{r(t)} p(t, w) dw (1 - k(0)) - L \right)^- \chi\{q = 0\} \end{aligned} \quad (3.11)$$

where  $r(t) = T + q(t - r(t))/L$ .  $k(t) = F(q(t))$  where  $F$  is continuous and when  $F(q(t)) = 1$  (i.e. when  $q(t) = q_{max}$ ),  $k(t)$  is determined by

$$\int_w \frac{w}{r(t)} p(t, w) dw \cdot (1 - k(t)) = L.$$

**Proof** If the initial distribution  $\mu$  has a continuous density then  $\mu(t, dw)$  has a continuous density  $p(t, w)$  differentiable in  $t$ . Hence (3.8) and (3.9) become

$$\begin{aligned} & \int_w (p(t, w) - p(0, w)) g(w) dw \\ &= \int_{s=0}^t \int_w \int_v g(w) \mu(s - r(s), dv; s, d(2w)) \frac{v}{r(s)} k(s - r(s)) ds \\ & - \int_{s=0}^t \int_w \int_v g(w) \mu(s - r(s), dv; s, dw) \frac{v}{r(s)} k(s - r(s)) ds \\ & + \int_{s=0}^t \int_w \frac{1}{r(s)} (1 - k(s - r(s))) \frac{dg}{dw} \mu(s, dw) ds \\ &= \int_w g(w) \left( \int_{s=0}^t \left( (4wp(s, 2w) - wp(s, 2w)) \frac{k(s - r(s))}{r(s)} dw ds \right. \right. \\ & \left. \left. - \int_w g(w) \left( \int_{s=0}^t \frac{1}{r(s)} (1 - k(s - r(s))) \frac{\partial p(s, w)}{\partial w} \right) dw \right) \right) \end{aligned}$$

and

$$\begin{aligned} q(t) - q(0) &= \int_{s=0}^t \left( \int_w p(s, w) \frac{w}{r(s)} (1 - k(s)) - L \right. \\ & \left. + \left( \int_w p(s, w) \frac{w}{r(s)} (1 - k(0)) - L \right)^- \chi\{q(t) = 0\} \right) ds. \end{aligned}$$

Since  $g$  is arbitrary we have the Fokker-Planck equation of the theorem. This system has a unique solution because it has no singularity. ■

An even grosser approximation is obtained by taking  $g(x) = x$  in equation (3.7). Assuming  $w(t) = \langle w, \mu(t, dw) \rangle$  is finite we obtain

$$\begin{aligned} w(t) - w(0) &= \int_0^t \left[ \frac{1}{r(s)} (1 - k(s - r(s))) \right. \\ & + \left( \langle wv, \mu(s - r(s), dv; s, d(2w)) \rangle - \langle wv, \mu(s - r(s), dv; s, d(w)) \rangle \right) \frac{1}{r(s)} k(s - r(s)) \Big] ds \\ &= \int_0^t \left[ \frac{1}{r(s)} (1 - k(s - r(s))) - \frac{1}{2} \langle wv, \mu(s - r(s), dv; s, d(w)) \rangle \frac{1}{r(s)} k(s - r(s)) \right] ds \quad (3.12) \end{aligned}$$

If we make the approximation that the window size at time  $s$  is uncorrelated with the window size one RTT ago then the second term in (3.12) becomes

$$\frac{1}{2} w(s - r(s)) w(s) \frac{1}{r(s)} k(s - r(s)).$$

This is the term in Equation (1) in [9]. Note that if the loss rate is small then the above term is small so the results can still be good. Equation (1) in [10] has a term like this as well (but the covariance between  $W(t)$  and  $W(t - r(t))$  is approximated by a variance).

A more precise approximation can be obtained by investigating the evolution of the window size  $W(t)$  of a canonical connection having distribution  $\mu(t; dw)$ . As a rough approximation there are two ways of arriving at a window size of  $w$  at time  $t$ . It happens if  $W(t - r(t)) = w - 1$  and there are no losses in the round trip time before  $t - r(t)$ . If the  $w - 1$  packets were evenly distributed across that RTT then the approximate probability there were no losses is approximately  $H(t - r(t), w - 1)$  where

$$H(t, w) = \prod_{j=1}^w \left( 1 - k(t - \frac{r(t)}{w}j) \right).$$

It also happens if  $W(t - r(t)) = 2w$  and there is one losses in the round trip time before  $t - r(t)$ . If the  $2w$  packets were evenly distributed across that RTT then the approximate probability there was at least one loss is approximately  $1 - H(t - r(t), 2w)$ . In the second case, we will neglect events where a loss just before time  $t - r(t)$  plus a loss at time  $t - r(t)$  causes a timeout. We also will neglect the possibility the window was  $4w$  between two and three round trip times ago and so on.

Using time reversal we have,

$$\langle v, \mu(s - r(s), dv; s, dw) \rangle \quad (3.13)$$

$$\begin{aligned} &= (w - 1)P(W(t) = w | W(t - r(t)) = w - 1)\mu(t - r(t), dw - 1) \\ &\quad + (2w)P(W(t) = w | W(t - r(t)) = 2w)\mu(t - r(t), d(2w)) \\ &= (w - 1)H(t - r(t), w - 1)\mu(t - r(t), dw - 1) \\ &\quad + (2w)(1 - H(t - r(t), 2w))\mu(t - r(t), d(2w)) \quad (3.14) \\ &= (w - 1)H(t - r(t), w - 1)p(t - r(t), w - 1)dw \\ &\quad + (2w)(1 - H(t - r(t), 2w))p(t - r(t), 2w)2dw \\ &= \left[ (w - 1)H(t - r(t), w - 1) \frac{p(t - r(t), w - 1)}{p(t, w)} \right. \\ &\quad \left. + (2w)(1 - H(t - r(t), 2w)) \frac{2p(t - r(t), 2w)}{p(t, w)} \right] p(t, w)dw \\ &= e(t; w)p(t, w)dw \end{aligned}$$

where  $e(t, w)$  is

$$(w - 1)H(t - r(t), w - 1) \frac{p(t - r(t), w - 1)}{p(t, w)} + (2w)(1 - H(t - r(t), 2w)) \frac{2p(t - r(t), 2w)}{p(t, w)}$$

Hence, by the same argument as Corollary 1 we get the refinement:

$$\begin{aligned} \frac{\partial p(t, w)}{\partial t} &= \left( p(t, 2w) \frac{e(t; 2w)}{r(t)} 2 - p(t, w) \frac{e(t; w)}{r(t)} \right) k(t - r(t)) \\ &\quad - \frac{1}{r(t)} (1 - k(t - r(t))) \frac{\partial p(t, w)}{\partial w} \quad (3.15) \end{aligned}$$

where  $r(t) = T + q(t - r(t))/L$  and  $q(t)$  satisfies (3.11).

From (3.14) we also get

$$\begin{aligned}
& \langle vw, \mu(s - r(s), dv; s, dw) \rangle \\
&= \langle w, (w - 1)H(t - r(t), w - 1)\mu(t - r(t), dw - 1) \rangle \\
&\quad + \langle w, (2w)(1 - H(t - r(t), 2w))\mu(t - r(t), d(2w)) \rangle \\
&= \langle (w + 1), wH(t - r(t), w)\mu(t - r(t), dw) \rangle + \frac{1}{2} \langle w, w(1 - H(t - r(t), w))\mu(t - r(t), d(w)) \rangle \\
&= \langle (w + 1)wH(t - r(t), w) + \frac{1}{2}w^2(1 - H(t - r(t), w)), \mu(t - r(t), dw) \rangle.
\end{aligned}$$

With this evaluation the second term in (3.12) becomes

$$-\frac{1}{2} \langle \frac{1}{2}w^2 + (\frac{w^2}{2} + w)H(t - r(t), w), \mu(t - r(t), dw) \rangle \frac{1}{r(t)} k(t - r(t))$$

which is close to  $-\frac{k(t-r(t))}{2r(t)} E(W(t - r(t))^2 + W(t - r(t)))$  if  $k(t - r(t))$  is small.

## 4 Fixed points of the mean-field equations

When the RTT is sufficiently small the approximating system (3.10) and (3.11) stabilizes; that is  $q(t)$  tends to a constant  $q$  and consequently the RTT,  $r(t)$ , and the loss rate,  $F(q(t))$ , tend to constants  $r$  and  $k$ . As the RTT increases however the delayed feedback will start amplifying any perturbation from equilibrium as was discussed in [9]. Since our window histogram is centered about its mean it follows that the stability analysis in [9] (modulo the refinements we have suggested above) should predict the bifurcation point.

For stable systems (3.10) and (3.11) become:

$$(1 - k) \frac{df_k(w)}{dw} = k(2(2w)f_k(2w) - wf_k(w)) \quad (4.16)$$

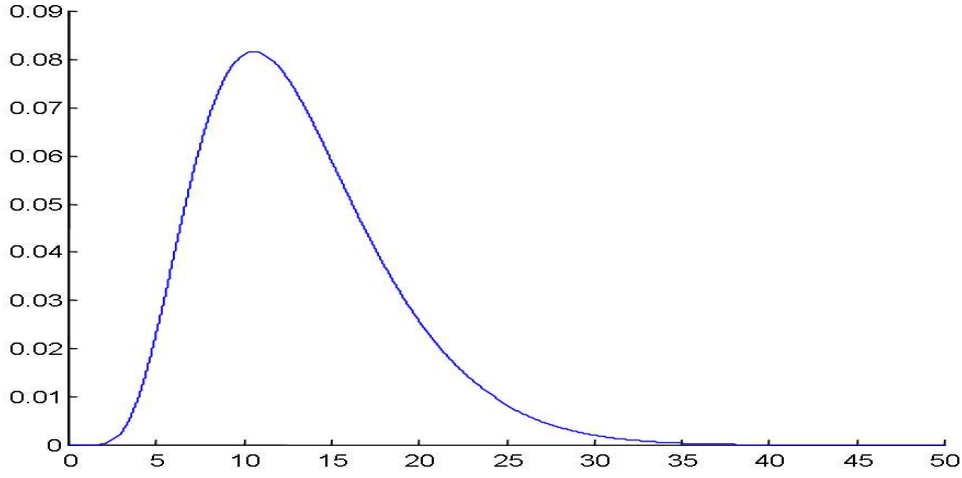
$$L = (1 - k) \frac{1}{r} \int_w wf_k(w) dw. \quad (4.17)$$

(4.17) is simply Little's formula since the right hand side represents the throughput as the average window size divided by the RTT times the proportion of packets that are not killed.

**Theorem 2** Let  $\Psi = \sum_{i=0}^{\infty} \frac{2^i}{\prod_{j=1}^i (1 - 4^j)}$  ( $\Psi \approx 0.4194$ ). The unique density  $f_k(w)$  solving (4.16) is given by

$$f_k(w) = \sum_{i=0}^{\infty} a_i \exp(-\frac{k}{1 - k} 4^i \frac{w^2}{2}) \quad (4.18)$$

$$a_0 = \sqrt{\frac{2}{\pi}} \frac{1}{\Psi} \sqrt{\frac{k}{1 - k}}; \quad a_i = a_{i-1} \frac{4}{1 - 4^i} = a_0 \frac{4^i}{\prod_{j=1}^i (1 - 4^j)}. \quad (4.19)$$

Figure 1: The histogram of window sizes in steady state when  $k = .01$ 

This solution was obtained in the paper by Adjih, Jacquet and Vvedenskaya [1] and independently by the authors.

**Proof** The fact that  $f_k$  is a solution follows by inspection. There are no convergence problems because the sum of the supremum norms of each term in (4.18) converges. The only point to elucidate is why  $f_k$  is a positive function : depending on the value of  $w$ , the modulus of the general term in the series decreases from the second or the first term on. Since the sign alternates,  $f_k(w) \geq \min(f_k^2(w), f_k^4(w)) \geq 0$ , where  $f_k^i(w)$  is the series trunked at the  $i^{th}$  term.

The value of  $a_0$  is determined by the requirement that  $f_k$  be a density.

$$\begin{aligned}
 1 &= \int_0^\infty f_k(w) dw = \sum_{i=0}^\infty a_i \int_0^\infty e^{-\frac{k}{1-k} \frac{w^2}{2} 4^i} dw = \sum_{i=0}^\infty a_i \frac{1}{2} \sqrt{2\pi \frac{1-k}{k 4^i}} \\
 &= \frac{\sqrt{2\pi}}{2} \sum_{i=0}^\infty a_i \frac{4^i}{\prod_{j=1}^i (1-4^j)} \sqrt{\frac{1-k}{k 4^i}} = \frac{\sqrt{1-k}}{\sqrt{k}} a_0 \sqrt{\frac{\pi}{2}} \sum_{i=0}^\infty \frac{2^i}{\prod_{j=1}^i (1-4^j)}
 \end{aligned} \tag{4.20}$$

■

There always exists a unique solution to both (4.16) and (4.17). First note that

$$\begin{aligned}
 \int_w w f_k(w) dw &= \sum_{i=0}^\infty a_i \int_0^\infty w \exp\left(-\frac{k}{1-k} 4^i \frac{w^2}{2}\right) dw \\
 &= \sum_{i=0}^\infty a_i \frac{1-k}{k 4^i} \\
 &= \frac{1-k}{k} a_0 \sum_{i=0}^\infty \frac{1}{\prod_{j=1}^i (1-4^j)} = \frac{1-k}{k} a_0 \xi \\
 &= \alpha \sqrt{\frac{1-k}{k}}
 \end{aligned} \tag{4.21}$$

where  $\xi = \sum_{i=0}^\infty \frac{1}{\prod_{j=1}^i (1-4^j)}$  ( $\xi \approx 0.6885$ ) and  $\alpha = \sqrt{\frac{2}{\pi} \frac{\xi}{\Psi}}$  ( $\alpha \approx 1.310$ ).



It follows from (4.17) that

$$L = \frac{(1-k)}{r} \alpha \sqrt{\frac{1-k}{k}} = \frac{\alpha}{r} \frac{(1-k)^{3/2}}{\sqrt{k}}.$$

Hence we get

$$\frac{(1-k)^3}{k} = \left( \frac{rL}{\alpha} \right)^2. \quad (4.22)$$

The function  $(1-k)^3/k$  is monotonically decreasing. It follows that there is a unique point  $k$  satisfying (4.22).

There are two possibilities. From (4.22), since  $(1-F(q))^3/F(q)$  is decreasing in  $q$ , the stable queue size is determined by the unique solution to

$$\frac{(1-F(q))^3}{F(q)} = \left( \frac{(T+q/L)L}{\alpha} \right)^2. \quad (4.23)$$

If the solution to (4.23) is less than  $q_{max}$  then  $k = F(q)$  gives the stable point and  $r = T + q/L$ . On the other hand if the solution to (4.23) is equal or greater than  $q_{max}$  then  $q = q_{max}$  and  $k$  is given by (4.22) (this is the case when  $Q_N(t)$  jitters at  $q_{max}$ ).

Until now we have ignored the fact that in equilibrium active connections go into timeout while an equal number become active (the slow-start period is assumed to be part of the timeout period so connections immediately enter congestion avoidance when they become active). At any time,  $N$  connections are active out of a total of  $N'$  connections and  $N' - N$ , the number of connections in timeout, is a proportion  $TO(k)$  of  $N$ . Hence  $N' - N = TO(k)N$ . In most cases we are given  $N'$  so solving the above equation for  $N$  gives the number of active connections.

The equilibrium window distribution can be used to calculate the number of connections in timeout. The long run proportion of connections in timeout,  $TO(k)$ , is equal to the proportion  $T(k)$  that enter timeout during one RTT times  $RTO/RTT$  where  $RTO$  is the timeout period (we take  $RTO$  to be one second). We could but will not analyse Selective Acknowledgements (SACK) or NewReno ([7]). We assume timeouts occur because there is a loss when the window size is less than four (we neglect the possibility of packets out of order). Since the number of losses in a window of size  $w$  has a Binomial distribution we have

$$\begin{aligned} T(k) &= \int_0^4 f_k(w)[1 - (1-k)^w]dw = \sum_{i=0}^{\infty} a_i T_i \quad \text{where} \\ T_i &= \frac{\sqrt{\pi}}{2} A_i \left[ erf\left(\frac{4}{A_i}\right) - \exp\left(\frac{A_i^2}{4} \ln(1-k)^2\right) \right. \\ &\quad \left. \left( erf\left(\frac{4}{A_i} - \ln(1-k)\frac{A_i}{2}\right) - erf\left(-\ln(1-k)\frac{A_i}{2}\right) \right) \right] \end{aligned} \quad (4.24)$$

where  $A_i = \sqrt{2(1-k)/(4^i k)}$  and  $erf(z) = (2/\sqrt{\pi}) \int_0^z \exp(-t^2)dt$ . The proof is given in the Appendix along with the Laplace transform of the window distribution.

To return from timeout a connection goes through slow-start until congestion avoidance starts after a packet loss. If the loss rate is  $k$  then the number of packets through the link from this slow-start is  $(1-k)/k$  on average, the mean of a geometric. These must be accounted for by modifying (4.17). Over time  $T$ , the amount of time in timeout is  $T \cdot N \cdot TO(k)$  and since each timeout lasts  $RTO$  this

means  $T \cdot N \cdot TO(k)/RTO$  timeouts are generated each of which generates  $(1 - k)/k$  packets on average; i.e.  $NTO(k)/RTO(1 - k)/k$  packets per second on average are generated by slow-start. Matching the link rate of the router with the incoming rate (4.17) becomes

$$NL = N \frac{(1 - k)}{r} \alpha \sqrt{\frac{1 - k}{k}} + N \frac{TO(k)}{RTO} \frac{1 - k}{k} \quad (4.25)$$

where  $r = T + q/L$ . We must now recalculate (4.22) by substituting  $k = F(q)$  into (4.25) thus determining  $k$ ,  $q$  and  $r$ .

**Theorem 3** *The number of active connections among  $N'$  connections is  $N$  where  $N' - N = TO(k)N$  and  $TO(k) = T(k)RTO/r$  where  $RTO$  is the timeout period plus the slowstart period,  $T(k)$  is given by (4.24) and  $k$  (and hence  $q$  and  $r$ ) is determined by (4.25).*

In addition, the equilibrium distribution of the window sizes (4.18) provides interesting QoS predictions. In particular the  $p^{th}$  quantile of (4.18)  $w_p$  provides that a connection spends a proportion  $p$  of its time with a window size less than  $w_p$ .

## 5 Analysis of the mean-field system

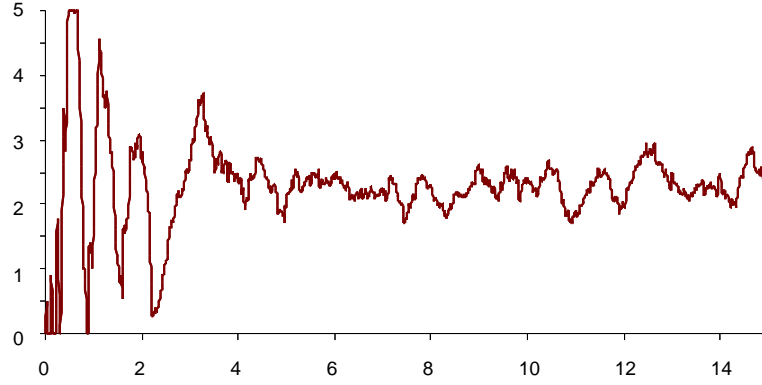
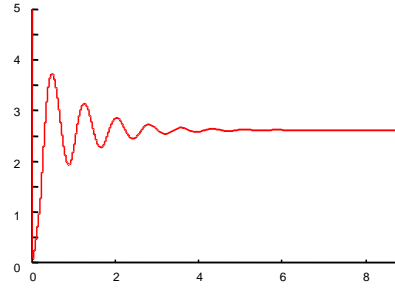
The system (3.15) and (3.11) can be solved numerically and it displays a whole range of behavior depending on the function  $F$  and the parameters  $L$  and  $T$ . To simplify matters we will take  $F(q) = 0$  if  $q < q_{min}$  and  $F(q) = 1$  if  $q \geq q_{max}$ . This covers the usual drop-tail case if  $q_{max} = B$ . We first look for fixed points as discussed in Section 4. Then we discuss a case where limit cycles are present and finally we discuss a case where there are multiple fixed points.

### 5.1 Numerical Results

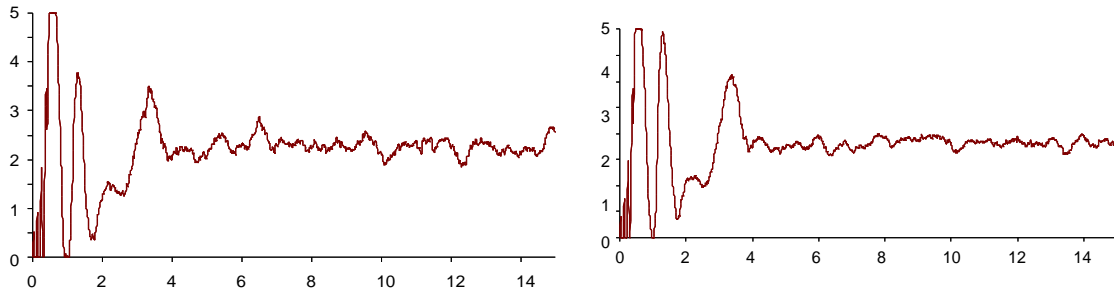
In the following numerical example  $N' = 200$  and the link rate is 44.736 Megabits per second. We assume each packet is 538 bytes to the link rate 10433 packets per second. We first assume there are no connections in timeout so  $N = N'$  so with  $N = 200$ ,  $L = 20.866$  packets per link per second. The transmission delay is  $T = .1$  seconds.  $Q_{min} = 0$ ,  $p_{min} = 0$ ,  $Q_{max} = 1000$ ,  $p_{max} = .05$ . By (4.23), the stable loss rate is  $k = 2.6\%$  at stable queue size  $q = 520.9$  and the RTT is 0.15 seconds. For the same parameters the numerical solution to the approximation (3.10) and (3.11) and to the approximation (3.15) and (3.11) both stabilize at the same queue length, loss rate and RTT. This is not surprising since the loss rate is quite low.

The Opnet simulation with 200 sources gives an average queue size of 452, an average loss rate of about 2.5% and 0.145 is the average RTT. Opnets lower loss rate combined with *lower* mean queue size results from timeouts which occur often because the mean window size is small. If we account for timeouts then solving  $200 - N = TO(N)N$  gives  $N = 192$  with 8 connections in timeout or slowstart. By solving (4.25)  $k = 0.0256$  which determines the average queue size of 512 and  $r = .149$ .  $T(0.0238) = .006$  and  $TO = 0.041$ . With the correction for timeouts we see the Opnet values are predicted slightly better. Timeouts may occur when there are multiple losses in a large window and this has not been modelled. The calculations given here should be refined.

The time series of queue sizes of the Opnet simulation is given in Figure 2 while the numerical solution using Matlab gives Figure 3.

Figure 2: Opnet simulation with  $N=200$ : the queue size divided by  $N$ Figure 3: Matlab solution with  $N=186$ : the queue size divided by  $N$ 

Next we redo the simulation with the same parameters but with  $N = 400$  connections and  $N = 800$  connections (see Figure 4).

Figure 4: Opnet simulation with  $N=400$  (on the left)  $N=800$  (on the right): the queue size divided by  $N$ 

We remark that the same behaviour occurs but the standard deviation of the oscillations around the average relative queue size gets smaller (probably proportionate to  $1/\sqrt{N}$ ).

Next we use the same parameters except that we increase the transmission delay is  $T = .3$  seconds. On the right side of Figure 5 the approximation (3.10) and (3.11) is indicated by a dotted line while the approximation (3.15) and (3.11) is indicated by a solid line. The approximation in [9] is given by the dashed line. All plots seem to oscillate around the steady state values but with bigger and bigger

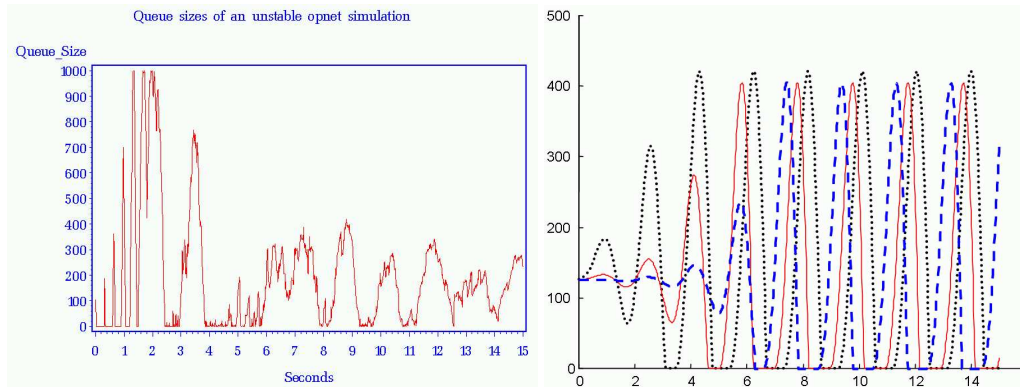


Figure 5: Queue sizes: Opnet on the left, Matlab on the right

oscillations until the queue hits zero. The Opnet simulation with 200 sources given on the left of Figure 5 is also unstable. We haven't corrected for proportion in timeout (which we can't calculate). Notice the period and amplitude of the oscillations observed in the Opnet simulation is pretty well predicted by all the approximations.

## 6 Conclusion

The mean-field model for the congestion windows of  $N$  TCP/IP sources multiplexed through a buffer implementing RED is given by (3.15) and (3.11). Simulation studies show excellent agreement when the number of sources is greater than 25. The equations for the evolution of the histogram of the window sizes provides a description of the quality of service experienced by each connection. The standard deviation of the histogram gives the variability of the throughput of a single connection and should be kept as small as possible. We can identify when the system becomes stable and when it becomes unstable because of the RTT delay in the feedback control loop. We can also identify systems with multiple equilibria caused by the RTT delay.

There are a host of outstanding questions associated with the analysis of the mean-field system. The most pressing is a derivation along the lines of [9] of the bifurcation point when a system goes from a stable queue size and loss rate to a system with limit cycles. It should also be possible to incorporate varying transmission time along with the window size of each connection to make a two dimensional histogram which along with the queue size more precisely describes the system. The mean field limit poses no additional difficulties.

## Acknowledgements

We thank Mike Maskery for his insights and his help with Matlab and Opnet. We also thank Michel Ouellette and Alan Chapman from Nortel Networks for their insight early on in this project.

## References

- [1] ADJIH, C., JACQUET, P., VVEDENSKAYA, N. (2001). Performance evaluation of a single queue under multi-user TCP/IP connections. *INRIA Research report #4141*.
- [2] AWEYA, J., OUELLETTE, M., DELFIN, Y. M., CHAPMAN, A. (2000). A load adaptive mechanism for buffer management. Nortel Networks Internal Report.

- [3] BRÉMAUD, P. (1981). Point Processes and Queues: Martingale Dynamics. *Springer Verlag*, 354 pp.
- [4] DAWSON, D. A. (1983). Critical Dynamics and Fluctuations for a Mean-Field model of cooperative behavior. *J. Statistical Phys.*, **31**, 29-85.
- [5] DIEKMANN, O. (1986). The Cell Size Distribution and Semigroups of Linear Operators. *Lecture notes in biomathematics: The dynamics of physiologically structured populations*; Metz, J.A.J ed. Springer
- [6] FLOYD, S., JACOBSON, V. (1993). Random early detection gateways for congestion avoidance. *IEEE/ACM Trans. Networking.*, **11**, No.4 397-413.
- [7] FLOYD, S. (1999). The NewReno modification to TCP's fast recovery algorithm. *RFC 2582*.
- [8] GROMOLL, H. C., PUHA, A. L., WILLIAMS, R. J. (2001). The fluid limit of a heavily loaded processor sharing queue. *preprint*.
- [9] HOLLOT, C.V., MISRA, V., TOWSLEY, D., GONG, W-B. (2001). A control theoretic analysis of RED. *To appear IEEE INFOCOM 2001*, 10 pp.
- [10] KUUSELA, P., LASSILA, P., VIRTAMO, J., KEY, P. (2001). Modeling RED with idealized TCP sources. *9<sup>th</sup> IFIP Conference on performance modelling and evaluation of ATM and IP networks 2001, Budapest*.
- [11] TINNACORNSRISUPHAP P., MAKOWSKI, A. (2001). Queue dynamics of RED gateways under a large number of TCP flows. *Globecom 2001*.

## 7 Appendix

From (4.18),

$$T(k) = \sum_{i=0}^{\infty} a_i \int_0^4 [1 - \exp(w \ln(1-k) - \frac{w^2}{A_i^2})] dw = \sum_{i=0}^{\infty} a_i T_i$$

where

$$\begin{aligned} T_i &= \int_0^4 [\exp(-\frac{w^2}{A_i^2}) - \exp(w \ln(1-k) - \frac{w^2}{A_i^2})] dw \\ &= \int_0^4 \exp(-\frac{w^2}{A_i^2}) dw - \exp(\frac{A_i^2}{4} \ln^2(1-k)) \int_0^4 \exp(-(\frac{w}{A_i} - A_i \ln(1-k)/2)^2) dw. \end{aligned}$$

A change of variable gives (4.24).

$$\text{Let } g_\delta(x) = \exp(-\delta \frac{x^2}{2}) \text{ and let } \hat{g}(t) = \int_0^\infty g(x) e^{tx} dx.$$

We will evaluate the Laplace transform of  $f_k$ ,

$$T_f(\theta) = \int_0^\infty f_k(w) e^{-\theta w} dw = \sum_{i=0}^{\infty} a_i \hat{g}_{\frac{k}{1-k} 4^i}(-\theta) \text{ by (4.18).} \quad (7.26)$$

We next calculate  $\hat{g}_\delta$ .

$$\begin{aligned}\hat{g}_\delta'(t) &= \int_0^\infty x g_\delta(x) e^{tx} dx = \int_0^\infty \frac{-1}{\delta} \frac{d(e^{\delta \frac{-x^2}{2}})}{dx} e^{tx} dx \\ &= -\left[ \frac{1}{\delta} e^{\delta \frac{-x^2}{2}} e^{tx} \right]_{x=0}^{x=\infty} + \frac{t}{\delta} \int_0^\infty g_\delta(x) e^{tx} dx \\ &= \frac{1}{\delta} + \frac{t}{\delta} \hat{g}_\delta(t), \text{ for all } t \in \mathbb{R}\end{aligned}$$

We can solve this equation by multiplying both sides by  $e^{-t^2/(2\delta)}$  which gives

$$\hat{g}_\delta'(t) e^{-t^2/(2\delta)} - \left( \frac{t}{\delta} e^{-t^2/(2\delta)} \right) \hat{g}_\delta(t) = \frac{1}{\delta} e^{-t^2/(2\delta)}.$$

Since  $\hat{g}_\delta(0) = 1/\sqrt{\pi/(2\delta)}$ ,

$$\hat{g}_\delta(t) = \sqrt{\frac{\pi}{2\delta}} + \frac{1}{\delta} e^{t^2/(2\delta)} \int_0^t e^{-x^2/(2\delta)} dx.$$

If we now substitute into (7.26) we get

$$T_f(\theta) = 1 - \sum_{i=0}^{\infty} a_i \frac{1-k}{4^i k} e^{\frac{\theta^2}{2} \frac{1-k}{k 4^i}} \int_0^\theta e^{-\frac{x^2}{2} \frac{1-k}{k 4^i}} dx$$

using (4.20).



---

Unité de recherche INRIA Rocquencourt  
Domaine de Voluceau - Rocquencourt - BP 105 - 78153 Le Chesnay Cedex (France)

Unité de recherche INRIA Lorraine : LORIA, Technopôle de Nancy-Brabois - Campus scientifique  
615, rue du Jardin Botanique - BP 101 - 54602 Villers-lès-Nancy Cedex (France)

Unité de recherche INRIA Rennes : IRISA, Campus universitaire de Beaulieu - 35042 Rennes Cedex (France)

Unité de recherche INRIA Rhône-Alpes : 655, avenue de l'Europe - 38330 Montbonnot-St-Martin (France)

Unité de recherche INRIA Sophia Antipolis : 2004, route des Lucioles - BP 93 - 06902 Sophia Antipolis Cedex (France)

---

Éditeur  
INRIA - Domaine de Voluceau - Rocquencourt, BP 105 - 78153 Le Chesnay Cedex (France)  
<http://www.inria.fr>  
ISSN 0249-6399